

به نام خدا

تمرین عملی شماره یک

مبانی بازیابی اطلاعات و جستجو وب

۹۸۱۶۰۷۳

دانیال بیاتی

۹۸۲۳۷۹۳

محمد جواد کفایتی

استاد : دکتر هدی مشایخی

برنامه ای بنویسید که تعدادی از اسناد به دلخواه دریافت کند و پرس و جو هایی را از کاربر دریافت کرده شامل **AND** و **OR** و **NOT** و به روش ایندکس معکوس اسناد مرتبط با آن را بازگرداند . اسناد را میتوانید از مجموعه داده های متنی از مخزن های متداول وب مثل **Kaggle** و **UCI** برداشت کنید. تعداد سندها نیازی نیست که زیاد باشد، و انجام کارهای پیش پردازش در این بخش نیازی نیست .نتایج پرس و جو را تحلیل کرده و درستی آنها را بررسی کنید.

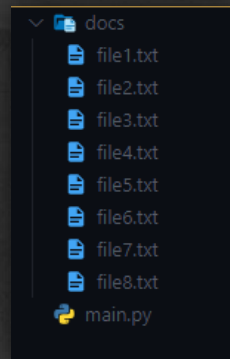
داخل پوشه **docs** حاوی ۸ داکيومنت است که به صورت **file1 ... file8** نام گذاری شده است که شامل متن هایی به صورت زیر است :

```
docs > file1.txt
```

```
1 I'm looking for information on the passing of author and screenwriter Nora Ephron.
```

```
docs > file5.txt
```

```
1 A friend is an identical twin. Her sister is already showing signs of dementia at 65.  
2 Now her twin is concerned she too will experience early onslaught of the disease.  
3 I would like to know more about the disease focused on females.  
4 Are there things we can do to lessen our chances of getting the disease?  
5 Is the disease always identified by genetic testing?
```



فایل **main.py** برنامه اصلی می باشد که از دو تابع ساخت ایندکس و سرچ کوئری قرار دارد.

تابع **create_inverted_index** :

به ترتیب ابتدا محتوای فایل های ورودی را خوانده و به کلمات را به صورت یک آرایه به نام **words** در می آوریم.

سپس اگر کلمه وجود نداشت در دیکشنری آن کلمه را اضافه می کنیم در غیر این صورت ایندکس فایل را که آن کلمه در آن موجود است را به کلمه مورد نظر در دیکشنری اضافه می کنیم این کار را تاجایی ادامه می دهیم که دیگر فایلی برای خواندن نداشته باشیم.

سپس هر کلمه در دیکشنری را با توجه به شماره فایل مرتب می کنیم



```
1 def create_inverted_index(files):
2     inverted_index = {}
3     for i, file in enumerate(files):
4         with open(file, 'r') as f:
5             contents = f.read()
6             words = contents.split()
7             for word in set(words):
8                 if word not in inverted_index:
9                     inverted_index[word] = []
10                    inverted_index[word].append(i + 1)
11     for word in inverted_index:
12         inverted_index[word].sort()
13     return inverted_index
14
15
16 files = glob.glob('./docs/*.txt')
17 inverted_index = create_inverted_index(files)
18 # print(inverted_index)
```



```
1 def search_query(inverted_index, query):
2     query_tokens = query.split()
3     result_docs = set(range(1, len(files) + 1))
4
5     # Loop through each token in the query
6     for i in range(len(query_tokens)):
7         token = query_tokens[i]
8
9         # If the token is "and"
10        if token == "and":
11            # Get the intersection of the current result document set and the next token's document set
12            next_token = query_tokens[i+1]
13            if next_token not in inverted_index:
14                return []
15            next_docs = set(inverted_index[next_token])
16            result_docs = result_docs.intersection(next_docs)
17
18        # If the token is "or"
19        elif token == "or":
20            # Get the union of the current result document set and the next token's document set
21            next_token = query_tokens[i+1]
22            if next_token not in inverted_index:
23                continue
24            next_docs = set(inverted_index[next_token])
25            print(result_docs)
26            result_docs = result_docs.union(next_docs)
27            print(result_docs)
28
29        # If the token is "not"
30        elif token == "not":
31            # Remove the documents from the current result document set that are also in the next token's document set
32            next_token = query_tokens[i+1]
33            if next_token not in inverted_index:
34                continue
35            next_docs = set(inverted_index[next_token])
36            result_docs = result_docs.difference(next_docs)
37
38        # If the token is not an operator, it must be a search term
39        else:
40            # Get the document set for the current token and update the result document set accordingly
41            if i - 1 >= 0 and (query_tokens[i - 1] == 'not' or query_tokens[i - 1] == 'or'):
42                continue
43            if token not in inverted_index:
44                return []
45            docs = set(inverted_index[token])
46            result_docs = result_docs.intersection(docs)
47
48        # Convert the final result document set to a sorted list and return it
49        result_list = sorted(list(result_docs))
50        return result_list
```

تابع `search_query` :

تابع `search_query` یک متن دریافت می‌کند و با استفاده از فهرست پوشا، اسنادی را که حاوی تمام کلمات موجود در متن هستند، پیدا می‌کند.

در ابتدا، متن ورودی به کلمات تقسیم شده و سپس برای هر کلمه، مجموعه اسنادی که آن کلمه در آن‌ها وجود دارد، با استفاده از فهرست پوشا به دست می‌آید. سپس با استفاده از عملگرهای منطقی `and` ، `or` و `not` ، مجموعه اسناد نهایی به دست می‌آید و به صورت لیست مرتب شده‌ای از اسناد بازگشت داده می‌شود.

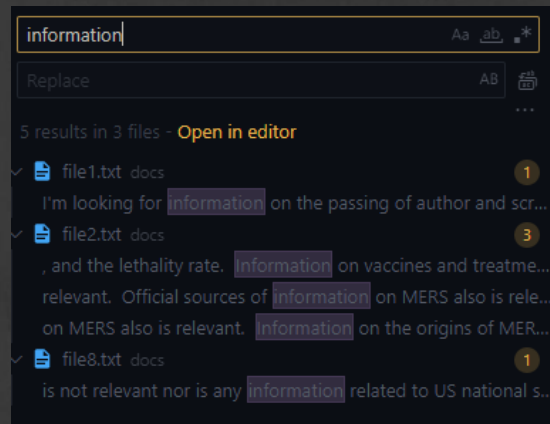
اگر کلمات جستجو شده در فهرست پوشا وجود نداشتند، خروجی خالی برگشت داده می‌شود.

```

1 files = glob.glob('./docs/*.txt')
2 inverted_index = create_inverted_index(files)
3 # print(inverted_index)
4
5 query = input("Please enter query : ")
6 print(search_query(inverted_index, query))

```

Please enter query : not information and virus
[3, 7]



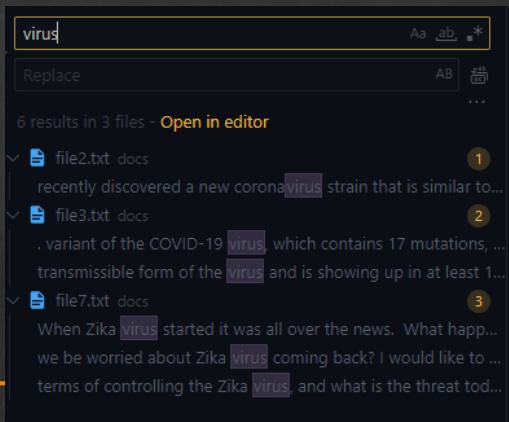
تست و بررسی داکيومنت ها :

برای مثال اگر کوئری زیر را داشته باشیم:

not information and virus

خروجی داکيومنت های ۳ و ۷ رانشان می دهد

همچنین با توجه به اینکه کلمه information در داکيومنت های ۱ و ۲ و ۸ وجود دارد لذا در داکيومنت ۳ و ۴ و ۵ و ۶ و ۷ وجود ندارد که اشتراکش با کلمه virus که در داکيومنت های ۳ و ۷ وجود دارد همان ۳ و ۷ می باشد.



فایل های متنی تمرین قبل را در تعداد بیشتر به وسیله نرم افزار **Lucene** ایندکس کنید. همه حروف را کوچک کنید و کلمات ایست را حذف کنید. ریشه یابی را به دلخواه خود می توانید انجام دهید. سپس حداقل ۵ پرس و جو مطرح کنید که عملگرهای **not, and, or** و جستجوی عبارت و پرس و جوی **wildcard** را شامل شوند و نتایج را بررسی کرده و گزارش کنید که آیا نتایج و رتبه بندی آنها منطقی هستند یا خیر.

داخل پوشه **docs** حاوی 30 داکيومنت است که به صورت **file1 ... file30** نام گذاری شده است که شامل متن هایی به صورت زیر است :

```
file21.txt M X
IR_T1_practical > code2 > docs > file21.txt
You, 1 second ago | 1 author (You)
1 Short ribs prosciutto picanha filet mignon sirloin drumstick chislic strip steak, beef meatloaf cow ribeye buffalo.
2 Tri-tip buffalo picanha meatball.
3 Venison doner beef ribs turkey, pancetta leberkas drumstick pastrami swine meatloaf corned beef t-bone.
4 Pork chop porchetta turkey landjaeger tongue shoulder leberkas drumstick.
5 Pork loin shankle tongue turducken shoulder jowl flank turkey pig.
6 Swine drumstick meatball, tongue pork chop kielbasa flank short loin frankfurter shank tenderloin meatloaf. You,
```

فایل **main.py** برنامه اصلی می باشد که از کتابخانه **whoosh** مثل **lucene** استفاده شده است که از توابع پیشفرض برای نرمالایز کردن داکيومنت ها استفاده شده است.

تابع **create_index_directory** : این تابع یک پوشه برای نگهداری ایندکس های ساخته شده می سازد.

تابع **create_schema** : این تابع از توابع پیشفرض در **whoosh** استفاده کرده و عملیاتی چون **stemming** و حذف **stopwords** ها را برعهده دارد که ایندکس های نهایی به صورت نرمالایز شده در پوشه ایندکس قرار می گیرد.

```
1 def create_index_directory(path="", name="index"):
2     if not os.path.exists(path+"index"):
3         os.mkdir(path+"index")
4     return create_in("index", schema)
5
6 def create_schema():
7     stopwords = frozenset(stoplists["en"])
8
9     analyzer = StemmingAnalyzer(stemfn=stem, stoplist=stopwords) | StopFilter(stoplist=stopwords)
10
11     return Schema(title=TEXT(stored=True), path=ID(stored=True), content=TEXT(analyzer=analyzer))
```

تابع **indexing_doc** : این تابع فایل های موجود در پوشه **docs** را به ترتیب ایندکس می کند و به عنوان یک **writer** فایل های ایندکس شده را برای جستجو و تحلیل **commit** می کند.

```
1 def indexing_docs(docs_addresses):
2     writer = ix.writer()
3
4     for i, file in enumerate(files):
5         with open(file, 'r') as f:
6             contents = f.read()
7
8             writer.add_document(title=file, path= file , content=contents)
9
10    writer.commit()
11
12
13 def searcher(boolean_query):
14     #for create parser and return docs number
15     with ix.searcher() as searcher:
16         parser = QueryParser("content", termclass=query.Variations ,schema=schema)
17         q = parser.parse(boolean_query)
18         results = searcher.search(q)
19         for i , r in enumerate(results):
20             title = r["title"]
21             print(f"score: {i+1}          -----          doc name: {title}")
22
```

تابع **searcher** : تابع **searcher** دارای پارامتر **boolean_query** می باشد که عبارت مورد نظر را دریافت می کند و به صورت پارس شده در جستجوگر قرار می دهد. این تابع با استفاده از **QueryParser** و **searcher** اطلاعات مورد نظر را به دست می آورد و شماره سند مورد نظر را بر می گرداند.



سپس با تست ۵ query صحت برنامه را امتحان می کنیم:



```
1 global schema , ix
2 files = glob.glob('./docs/*.txt')
3 schema = create_schema()
4 ix = create_index_directory()
5 indexing_docs(files)
6
7 query_input = input("Enter query : ")
8 searcher(query_input)
```

```
23:07:51 > py main.py
Enter query : starT OR Hello
score: 1      ----- doc name: ./docs\file24.txt
score: 2      ----- doc name: ./docs\file7.txt
score: 3      ----- doc name: ./docs\file5.txt
score: 4      ----- doc name: ./docs\file25.txt
score: 5      ----- doc name: ./docs\file16.txt
score: 6      ----- doc name: ./docs\file29.txt
score: 7      ----- doc name: ./docs\file21.txt
```

```
23:12:40 > py main.py
Enter query : HELLO AND reading
score: 1      ----- doc name: ./docs\file5.txt
```

```
23:13:30 > py main.py
Enter query : relevant AND NOT virus
score: 1      ----- doc name: ./docs\file2.txt
score: 2      ----- doc name: ./docs\file19.txt
score: 3      ----- doc name: ./docs\file18.txt
score: 4      ----- doc name: ./docs\file14.txt
score: 5      ----- doc name: ./docs\file8.txt
score: 6      ----- doc name: ./docs\file17.txt
score: 7      ----- doc name: ./docs\file11.txt
```

```
23:11:39 > py main.py
Enter query : HELLO AND NOT READING
score: 1      ----- doc name: ./docs\file25.txt
score: 2      ----- doc name: ./docs\file16.txt
score: 3      ----- doc name: ./docs\file29.txt
score: 4      ----- doc name: ./docs\file21.txt
```

```
23:10:27 > py main.py
Enter query : Now AND ReaD
score: 1      ----- doc name: ./docs\file17.txt
score: 2      ----- doc name: ./docs\file5.txt
```